**Data-Intensive Information Processing Applications — Session #12**

# Bigtable, Hive, and Pig

**Jimmy Lin**
University of Maryland

Tuesday, April 27, 2010

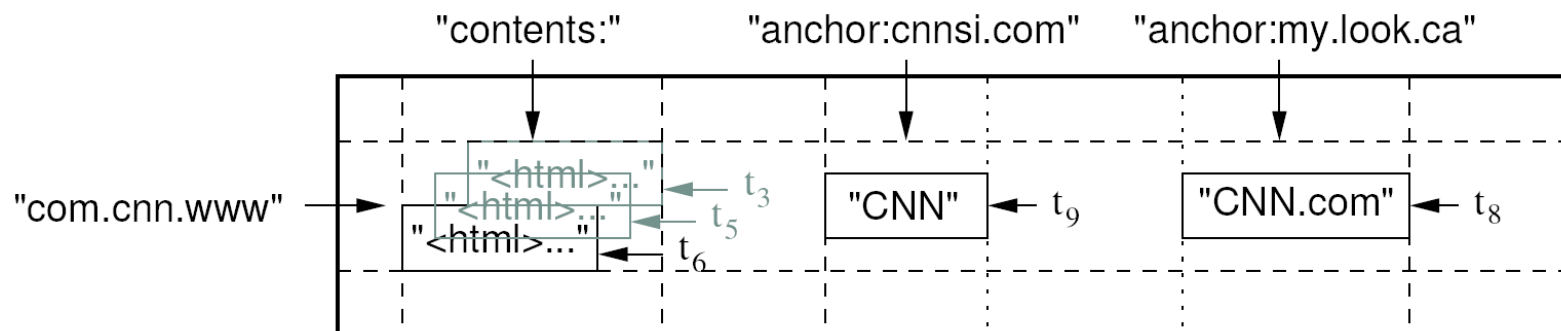Source: Wikipedia (Japanese rock garden)

# Today's Agenda

- Bigtable
- Hive
- Pig

# Bigtable

# Data Model

- A table in Bigtable is a sparse, distributed, persistent multidimensional sorted map

- Map indexed by a row key, column key, and a timestamp
  - (row:string, column:string, time:int64) $\rightarrow$ uninterpreted byte array

- Supports lookups, inserts, deletes
  - Single row transactions only

# Rows and Columns

- Rows maintained in sorted lexicographic order
  - Applications can exploit this property for efficient row scans
  - Row ranges dynamically partitioned into tablets

- Columns grouped into column families
  - Column key = *family:qualifier*
  - Column families provide locality hints
  - Unbounded number of columns

# Bigtable Building Blocks

- GFS

- Chubby

- SSTable

# SSTable

- Basic building block of Bigtable

- Persistent, ordered immutable map from keys to values
  - Stored in GFS

- Sequence of blocks on disk plus an index for block lookup
  - Can be completely mapped into memory

- Supported operations:
  - Look up value associated with key
  - Iterate key/value pairs within a key range

# Tablet

- Dynamically partitioned range of rows

- Built from multiple SSTables

| Tablet | Start:aardvark | | End:apple |
|---|---|---|---|

| 64K block | 64K block | 64K block | SSTable |
|---|---|---|---|
| | | | Index |

| 64K block | 64K block | 64K block | SSTable |
|---|---|---|---|
| | | | Index |

# Table

- Multiple tablets make up the table

- SSTables can be shared

```
┌─────────────────────────┐        ┌─────────────────────────┐
│ Tablet                  │        │ Tablet                  │
│                         │        │                         │
│ aardvark        apple   │        │ apple_two_E      boat   │
└─────────────────────────┘        └─────────────────────────┘
        │        │          ╲       ╱         │         │
        ▼        ▼           ╲     ╱          ▼         ▼
   ┌────────┐┌────────┐  ┌────────┐┌────────┐
   │SSTable ││SSTable │  │SSTable ││SSTable │
   │        ││        │  │        ││        │
   └────────┘└────────┘  └────────┘└────────┘
```

# Architecture

- Client library

- Single master server

- Tablet servers

# Bigtable Master

- Assigns tablets to tablet servers

- Detects addition and expiration of tablet servers

- Balances tablet server load

- Handles garbage collection

- Handles schema changes

# Bigtable Tablet Servers

- Each tablet server manages a set of tablets
  - Typically between ten to a thousand tablets
  - Each 100-200 MB by default

- Handles read and write requests to the tablets

- Splits tablets that have grown too large

# Tablet Location



Upon discovery, clients cache tablet locations

# Tablet Assignment

- Master keeps track of:

  - Set of live tablet servers
  - Assignment of tablets to tablet servers
  - Unassigned tablets

- Each tablet is assigned to one tablet server at a time

  - Tablet server maintains an exclusive lock on a file in Chubby
  - Master monitors tablet servers and handles assignment

- Changes to tablet structure

  - Table creation/deletion (master initiated)
  - Tablet merging (master initiated)
  - Tablet splitting (tablet server initiated)

# Tablet Serving



**"Log Structured Merge Trees"**

# Compactions

- Minor compaction

  - Converts the memtable into an SSTable
  - Reduces memory usage and log traffic on restart

- Merging compaction

  - Reads the contents of a few SSTables and the memtable, and writes out a new SSTable
  - Reduces number of SSTables

- Major compaction

  - Merging compaction that results in only one SSTable
  - No deletion records, only live data

# Bigtable Applications

- Data source and data sink for MapReduce

- Google's web crawl

- Google Earth

- Google Analytics

# Lessons Learned

- Fault tolerance is hard

- Don't add functionality before understanding its use
  - Single-row transactions appear to be sufficient

- Keep it simple!

# HBase

- Open-source clone of Bigtable

- Implementation hampered by lack of file append in HDFS

# Hive and Pig

# Need for High-Level Languages

○ Hadoop is great for large-data processing!

  ● But writing Java programs for everything is verbose and slow

  ● Not everyone wants to (or can) write Java code

○ Solution: develop higher-level data processing languages

  ● Hive: HQL is like SQL

  ● Pig: Pig Latin is a bit like Perl

# Hive and Pig

- Hive: data warehousing application in Hadoop
  - Query language is HQL, variant of SQL
  - Tables stored on HDFS as flat files
  - Developed by Facebook, now open source

- Pig: large-scale data processing system
  - Scripts are written in Pig Latin, a dataflow language
  - Developed by Yahoo!, now open source
  - Roughly 1/3 of all Yahoo! internal jobs

- Common idea:
  - Provide higher-level language to facilitate large-data processing
  - Higher-level language "compiles down" to Hadoop jobs

# Hive: Background

- Started at Facebook

- Data was collected by nightly cron jobs into Oracle DB

- "ETL" via hand-coded python

- Grew from 10s of GBs (2006) to 1 TB/day new data (2007), now 10x that

# Hive Components

- Shell: allows interactive queries

- Driver: session handles, fetch, execute

- Compiler: parse, plan, optimize

- Execution engine: DAG of stages (MR, HDFS, metadata)

- Metastore: schema, location in HDFS, SerDe

# Data Model

- Tables
  - Typed columns (int, float, string, boolean)
  - Also, list: map (for JSON-like data)

- Partitions
  - For example, range-partition tables by date

- Buckets
  - Hash partitions within ranges (useful for sampling, join optimization)

# Metastore

- Database: namespace containing a set of tables

- Holds table definitions (column types, physical layout)

- Holds partitioning information

- Can be stored in Derby, MySQL, and many other relational databases

# Physical Layout

- Warehouse directory in HDFS
  - E.g., /user/hive/warehouse

- Tables stored in subdirectories of warehouse
  - Partitions form subdirectories of tables

- Actual data stored in flat files
  - Control char-delimited text, or SequenceFiles
  - With custom SerDe, can use arbitrary format

# Hive: Example

- Hive looks similar to an SQL database

- Relational join on two tables:
  - Table of word counts from Shakespeare collection
  - Table of word counts from the bible

    SELECT s.word, s.freq, k.freq FROM shakespeare s
      JOIN bible k ON (s.word = k.word) WHERE s.freq >= 1 AND k.freq >= 1
      ORDER BY s.freq DESC LIMIT 10;

    | | | |
    |---|---|---|
    | the | 25848 | 62394 |
    | I | 23031 | 8854 |
    | and | 19671 | 38985 |
    | to | 18038 | 13526 |
    | of | 16700 | 34654 |
    | a | 14170 | 8057 |
    | you | 12702 | 2720 |
    | my | 11297 | 4135 |
    | in | 10797 | 12445 |
    | is | 8882 | 6884 |

# Hive: Behind the Scenes

SELECT s.word, s.freq, k.freq FROM shakespeare s
   JOIN bible k ON (s.word = k.word) WHERE s.freq >= 1 AND k.freq >= 1
   ORDER BY s.freq DESC LIMIT 10;

(Abstract Syntax Tree)

(TOK_QUERY (TOK_FROM (TOK_JOIN (TOK_TABREF shakespeare s) (TOK_TABREF bible k) (= (. (TOK_TABLE_OR_COL s) word) (. (TOK_TABLE_OR_COL k) word)))) (TOK_INSERT (TOK_DESTINATION (TOK_DIR TOK_TMP_FILE)) (TOK_SELECT (TOK_SELEXPR (. (TOK_TABLE_OR_COL s) word)) (TOK_SELEXPR (. (TOK_TABLE_OR_COL s) freq)) (TOK_SELEXPR (. (TOK_TABLE_OR_COL k) freq))) (TOK_WHERE (AND (>= (. (TOK_TABLE_OR_COL s) freq) 1) (>= (. (TOK_TABLE_OR_COL k) freq) 1))) (TOK_ORDERBY (TOK_TABSORTCOLNAMEDESC (. (TOK_TABLE_OR_COL s) freq))) (TOK_LIMIT 10)))

(one or more of MapReduce jobs)

# Hive: Behind the Scenes

```
STAGE DEPENDENCIES:
  Stage-1 is a root stage
  Stage-2 depends on stages: Stage-1
  Stage-0 is a root stage

STAGE PLANS:
  Stage: Stage-1
   Map Reduce
    Alias -> Map Operator Tree:
      s
       TableScan
        alias: s
        Filter Operator
         predicate:
           expr: (freq >= 1)
           type: boolean
         Reduce Output Operator
          key expressions:
            expr: word
            type: string
          sort order: +
          Map-reduce partition columns:
            expr: word
            type: string
          tag: 0
          value expressions:
            expr: freq
            type: int
            expr: word
            type: string
      k
       TableScan
        alias: k
        Filter Operator
         predicate:
           expr: (freq >= 1)
           type: boolean
         Reduce Output Operator
          key expressions:
            expr: word
            type: string
          sort order: +
          Map-reduce partition columns:
            expr: word
            type: string
          tag: 1
          value expressions:
            expr: freq
            type: int
```

```
Reduce Operator Tree:
    Join Operator
     condition map:
       Inner Join 0 to 1
     condition expressions:
       0 {VALUE._col0} {VALUE._col1}
       1 {VALUE._col0}
     outputColumnNames: _col0, _col1, _col2
     Filter Operator
      predicate:
        expr: ((_col0 >= 1) and (_col2 >= 1))
        type: boolean
      Select Operator
       expressions:
         expr: _col1
         type: string
         expr: _col0
         type: int
         expr: _col2
         type: int
       outputColumnNames: _col0, _col1, _col2
       File Output Operator
        compressed: false
        GlobalTableId: 0
        table:
          input format: org.apache.hadoop.mapred.SequenceFileInputFormat
          output format: org.apache.hadoop.hive.ql.io.HiveSequenceFileOutputFormat
```

```
Stage: Stage-2
 Map Reduce
  Alias -> Map Operator Tree:
    hdfs://localhost:8022/tmp/hive-training/364214370/10002
      Reduce Output Operator
       key expressions:
         expr: _col1
         type: int
       sort order: -
       tag: -1
       value expressions:
         expr: _col0
         type: string
         expr: _col1
         type: int
         expr: _col2
         type: int
  Reduce Operator Tree:
   Extract
    Limit
     File Output Operator
      compressed: false
      GlobalTableId: 0
      table:
        input format: org.apache.hadoop.mapred.TextInputFormat
        output format: org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputFormat


Stage: Stage-0
 Fetch Operator
  limit: 10
```

# Hive Demo

# Example Data Analysis Task

## Find users who tend to visit "good" pages.

**Visits**

| user | url | time |
|------|-----|------|
| Amy | www.cnn.com | 8:00 |
| Amy | www.crap.com | 8:05 |
| Amy | www.myblog.com | 10:00 |
| Amy | www.flickr.com | 10:05 |
| Fred | cnn.com/index.htm | 12:00 |

⋮

**Pages**

| url | pagerank |
|-----|----------|
| www.cnn.com | 0.9 |
| www.flickr.com | 0.9 |
| www.myblog.com | 0.7 |
| www.crap.com | 0.2 |

⋮

Pig Slides adapted from Olston et al.

# Conceptual Dataflow

Load
Visits(user, url, time)

Load
Pages(url, pagerank)

Canonicalize URLs

Join
url = url

Group by user

Compute Average Pagerank

Filter
avgPR > 0.5

# System-Level Dataflow

**Visits**  **Pages**

load

canonicalize

. . .  . . .  load

join by url

group by user

compute average pagerank

filter

the answer

# MapReduce Code

```java
import java.io.IOException;
import java.util.ArrayList;
import java.util.Iterator;
import java.util.List;

import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.io.Writable;
import org.apache.hadoop.io.WritableComparable;
import org.apache.hadoop.mapred.FileInputFormat;
import org.apache.hadoop.mapred.FileOutputFormat;
import org.apache.hadoop.mapred.JobConf;
import org.apache.hadoop.mapred.KeyValueTextInputFormat;
import org.a    pache.hadoop.mapred.Mapper;
import org.apache.hadoop.mapred.MapReduceBase;
import org.apache.hadoop.mapred.OutputCollector;
import org.apache.hadoop.mapred.RecordReader;
import org.apache.hadoop.mapred.Reducer;
import org.apache.hadoop.mapred.Reporter;
imp ort org.apache.hadoop.mapred.SequenceFileInputFormat;
import org.apache.hadoop.mapred.SequenceFileOutputFormat;
import org.apache.hadoop.mapred.TextInputFormat;
import org.apache.hadoop.mapred.jobcontrol.Job;
import org.apache.hadoop.mapred.jobcontrol.JobC            ontrol;
import org.apache.hadoop.mapred.lib.IdentityMapper;

public class MRExample {
    public static class LoadPages extends MapReduceBase
        implements Mapper<LongWritable, Text, Text, Text> {

        public void map(LongWritable k, Text val,
                OutputCollector<Text, Text> oc,
                Reporter reporter) throws IOException {
            // Pull the key out
            String line = val.toString();
            int firstComma = line.indexOf(',');
            String key = line.sub        string(0, firstComma);
            String value = line.substring(firstComma + 1);
            Text outKey = new Text(key);
            // Prepend an index to the value so we know which file
            // it came from.
            Text outVal = new Text("1             " + value);
            oc.collect(outKey, outVal);
        }
    }
    public static class LoadAndFilterUsers extends MapReduceBase
        implements Mapper<LongWritable, Text, Text, Text> {

        public void map(LongWritable k, Text val,
                OutputCollector<Text, Text> oc,
                Reporter reporter) throws IOException {
            // Pull the key out
            String line = val.toString();
            int firstComma = line.indexOf(',');
            String value = line.substring(                firstComma + 1);
            int age = Integer.parseInt(value);
            if (age < 18 || age > 25) return;
            String key = line.substring(0, firstComma);
            Text outKey = new Text(key);
            // Prepend an index to the value so w              e know which file
            // it came from.
            Text outVal = new Text("2" + value);
            oc.collect(outKey, outVal);
        }
    }
    public static class Join extends MapReduceBase
        implements Reducer<Text, Text, Text, Text> {

        public void reduce(Text key,
                Iterator<Text> iter,
                OutputCollector<Text, Text> oc,
                Reporter reporter) throws IOException {
            // For each value, figure out which file it's from and
store it
            // accordingly.
            List<String> first = new ArrayList<String>();
            List<String> second = new ArrayList<String>();

            while (iter.hasNext()) {
                Text t = iter.next();
                String value = t.to            String();
                if (value.charAt(0) == '1')
first.add(value.substring(1));
                else second.add(value.substring(1));
```

```java
                    reporter.setStatus("OK");
            }

            // Do the cross product and collect the values
            for (String s1 : first) {
                for (String s2 : second) {
                    String outval = key + "," + s1 + "," + s2;
                    oc.collect(null, new Text(outval));
                    reporter.setStatus("OK");
                }
            }
        }
    }
    public static class LoadJoined extends MapReduceBase
        implements Mapper<Text, Text, Text, LongWritable> {

        public void map(
                Text k,
                Text val,
                OutputColle         ctor<Text, LongWritable> oc,
                Reporter reporter) throws IOException {
            // Find the url
            String line = val.toString();
            int firstComma = line.indexOf(',');
            int secondComma = line.indexOf(',', first           Comma);
            String key = line.substring(firstComma, secondComma);
            // drop the rest of the record, I don't need it anymore,
            // just pass a 1 for the combiner/reducer to sum instead.
            Text outKey = new Text(key);
            oc.collect(outKey, new LongWritable(1L));
        }
    }
    public static class ReduceUrls extends MapReduceBase
        implements Reducer<Text, LongWritable, WritableComparable,
Writable> {

        public void reduce(
                Text ke          y,
                Iterator<LongWritable> iter,
                OutputCollector<WritableComparable, Writable> oc,
                Reporter reporter) throws IOException {
            // Add up all the values we see

            long sum = 0;
            wh      ile (iter.hasNext()) {
                sum += iter.next().get();
                reporter.setStatus("OK");
            }

            oc.collect(key, new LongWritable(sum));
        }
    }
    public static class LoadClicks extends MapReduceBase
        i   mplements Mapper<WritableComparable, Writable, LongWritable,
Text> {

        public void map(
                WritableComparable key,
                Writable val,
                OutputCollector<LongWritable, Text> oc,
                Reporter reporter)            throws IOException {
            oc.collect((LongWritable)val, (Text)key);
        }
    }
    public static class LimitClicks extends MapReduceBase
        implements Reducer<LongWritable, Text, LongWritable, Text> {

        int count = 0;
        public        void reduce(
            LongWritable key,
            Iterator<Text> iter,
            OutputCollector<LongWritable, Text> oc,
            Reporter reporter) throws IOException {

            // Only output the first 100 records
            while (count        < 100 && iter.hasNext()) {
                oc.collect(key, iter.next());
                count++;
            }
        }
    }
    public static void main(String[] args) throws IOException {
        JobConf lp = new JobConf(MRExample.class);
        lp.se    tJobName("Load Pages");
        lp.setInputFormat(TextInputFormat.class);
```

```java
        lp.setOutputKeyClass(Text.class);
        lp.setOutputValueClass(Text.class);
        lp.setMapperClass(LoadPages.class);
        FileInputFormat.addInputPath(lp, new
Path("/    user/gates/pages"));
        FileOutputFormat.setOutputPath(lp,
            new Path("/user/gates/tmp/indexed_pages"));
        lp.setNumReduceTasks(0);
        Job loadPages = new Job(lp);

        JobConf lfu = new JobConf(MRExample.class);
        lfu.s     etJobName("Load and Filter Users");
        lfu.setInputFormat(TextInputFormat.class);
        lfu.setOutputKeyClass(Text.class);
        lfu.setOutputValueClass(Text.class);
        lfu.setMapperClass(LoadAndFilterUsers.class);
        FileInputFormat.add        InputPath(lfu, new
Path("/user/gates/users"));
        FileOutputFormat.setOutputPath(lfu,
            new Path("/user/gates/tmp/filtered_users"));
        lfu.setNumReduceTasks(0);
        Job loadUsers = new Job(lfu);

        JobConf join = new JobConf(        MRExample.class);
        join.setJobName("Join Users and Pages");
        join.setInputFormat(KeyValueTextInputFormat.class);
        join.setOutputKeyClass(Text.class);
        join.setOutputValueClass(Text.class);
        join.setMapperClass(IdentityMap            per.class);
        join.setReducerClass(Join.class);
        FileInputFormat.addInputPath(join, new
Path("/user/gates/tmp/indexed_pages"));
        FileInputFormat.addInputPath(join, new
Path("/user/gates/tmp/filtered_users"));
        FileOutputFormat.se            tOutputPath(join, new
Path("/user/gates/tmp/joined"));
        join.setNumReduceTasks(50);
        Job joinJob = new Job(join);
        joinJob.addDependingJob(loadPages);
        joinJob.addDependingJob(loadUsers);

        JobConf group = new JobConf(MRE              xample.class);
        group.setJobName("Group URLs");
        group.setInputFormat(KeyValueTextInputFormat.class);
        group.setOutputKeyClass(Text.class);
        group.setOutputValueClass(LongWritable.class);
        group.setOutputFormat(SequenceFi            leOutputFormat.class);
        group.setMapperClass(LoadJoined.class);
        group.setCombinerClass(ReduceUrls.class);
        group.setReducerClass(ReduceUrls.class);
        FileInputFormat.addInputPath(group, new
Path("/user/gates/tmp/joined"));
        FileOutputFormat.setOutputPath(group, new
Path("/user/gates/tmp/grouped"));
        group.setNumReduceTasks(50);
        Job groupJob = new Job(group);
        groupJob.addDependingJob(joinJob);

        JobConf top100 = new JobConf(MRExample.class);
        top100.setJobName("Top 100 sites");
        top100.setInputFormat(SequenceFileInputFormat.class);
        top100.setOutputKeyClass(LongWritable.class);
        top100.setOutputValueClass(Text.class);
        top100.setOutputFormat(SequenceFileOutputF                ormat.class);
        top100.setMapperClass(LoadClicks.class);
        top100.setCombinerClass(LimitClicks.class);
        top100.setReducerClass(LimitClicks.class);
        FileInputFormat.addInputPath(top100, new
Path("/user/gates/tmp/grouped"));
        FileOutputFormat.setOutputPath(top100, new
Path("/user/gates/top100sitesforusers18to25"));
        top100.setNumReduceTasks(1);
        Job limit = new Job(top100);
        limit.addDependingJob(groupJob);

        JobControl jc = new JobControl("Find top             100 sites for users
18 to 25");
        jc.addJob(loadPages);
        jc.addJob(loadUsers);
        jc.addJob(joinJob);
        jc.addJob(groupJob);
        jc.addJob(limit);
        jc.run();
    }
}
```

Pig Slides adapted from Olston et al.

# Pig Latin Script

```
Visits = load    '/data/visits' as (user, url, time);
Visits = foreach Visits generate user, Canonicalize(url), time;

Pages = load    '/data/pages' as (url, pagerank);

VP = join    Visits by url, Pages by url;
UserVisits = group    VP by user;
UserPageranks = foreach UserVisits generate user,
AVG(VP.pagerank) as avgpr;
GoodUsers = filter  UserPageranks by avgpr > '0.5';

store    GoodUsers into '/data/good_users';
```
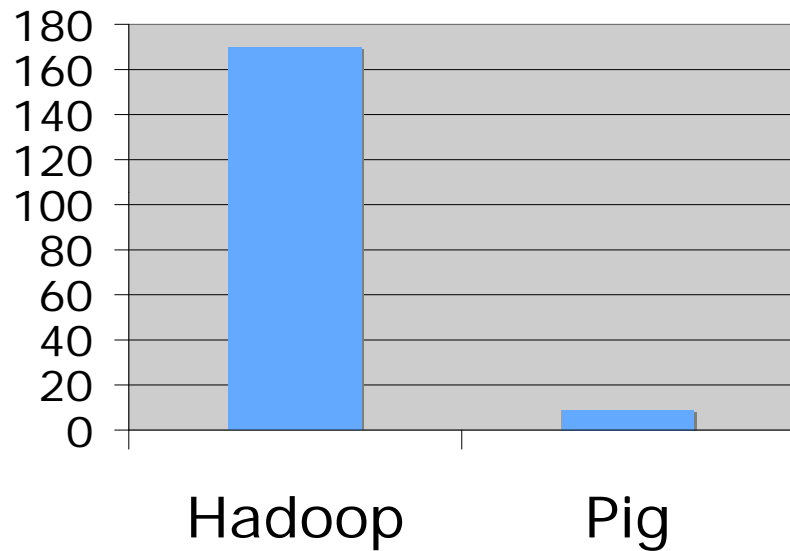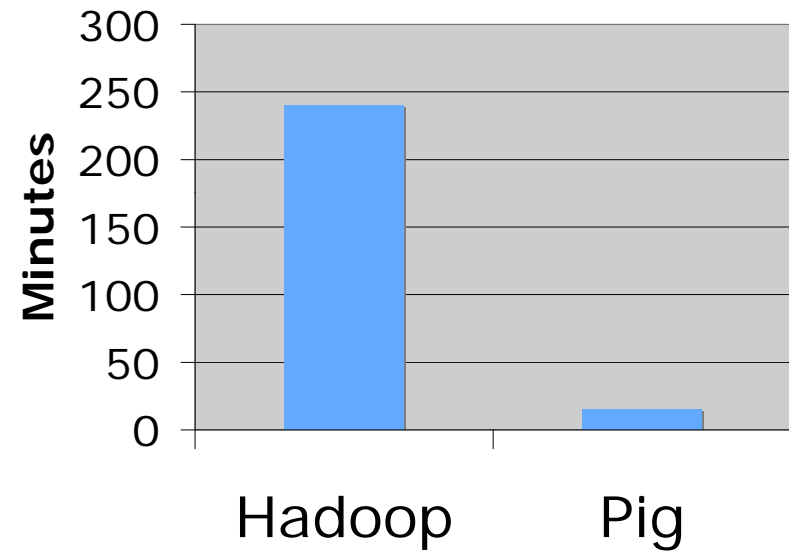
Pig Slides adapted from Olston et al.

# Java vs. Pig Latin

**1/20 the lines of code**



**1/16 the development time**



**Performance on par with raw Hadoop!**

Pig Slides adapted from Olston et al.

# Pig takes care of...

- Schema and type checking

- Translating into efficient physical dataflow

  - (i.e., sequence of one or more MapReduce jobs)

- Exploiting data reduction opportunities

  - (e.g., early partial aggregation via a combiner)

- Executing the system-level dataflow

  - (i.e., running the MapReduce jobs)

- Tracking progress, errors, etc.

# Pig Demo

Questions?

Source: Wikipedia (Japanese rock garden)